

```
In [1]: import nltk
        from nltk.tokenize import sent_tokenize, word_tokenize
```

```
In [6]: import nltk
        nltk.download('punkt')
```

```
[nltk_data] Downloading package punkt to
[nltk_data]   C:\Users\admin\AppData\Roaming\nltk_data...
[nltk_data]   Unzipping tokenizers\punkt.zip.
```

```
Out[6]: True
```

```
In [33]: example_string = """
... Muad'Dib learned rapidly because his first training was in how to learn.
... And the first lesson of all was the basic trust that he could learn.
... It's shocking to find how many people do not believe they can learn,
... and how many more believe learning to be difficult."""
```

```
In [34]: sent_tokenize(example_string)
```

```
Out[34]: ["\nMuad'Dib learned rapidly because his first training was in how to learn.",
          'And the first lesson of all was the basic trust that he could learn.',
          "It's shocking to find how many people do not believe they can learn,\nand how ma
ny more believe learning to be difficult."]
```

```
In [35]: word_tokenize(example_string)
```

```
Out[35]: ["Muad'Dib",
          'learned',
          'rapidly',
          'because',
          'his',
          'first',
          'training',
          'was',
          'in',
          'how',
          'to',
          'learn',
          '.',
          'And',
          'the',
          'first',
          'lesson',
          'of',
          'all',
          'was',
          'the',
          'basic',
          'trust',
          'that',
          'he',
          'could',
          'learn',
          '.',
          'It',
          "'s",
          'shocking',
          'to',
          'find',
          'how',
          'many',
          'people',
          'do',
          'not',
          'believe',
          'they',
          'can',
          'learn',
          ',',
          'and',
          'how',
          'many',
          'more',
          'believe',
          'learning',
          'to',
          'be',
          'difficult',
          '.']
```

```
In [9]: nltk.download("stopwords")
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
```

```
[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\admin\AppData\Roaming\nltk_data...
[nltk_data] Unzipping corpora\stopwords.zip.
```

```
In [ ]:
```

```
In [36]: worf_quote = "Sir, I protest. I am not a merry man!"
```

```
In [37]: words_in_quote = word_tokenize(worf_quote)
words_in_quote
```

```
Out[37]: ['Sir', ',', 'I', 'protest', '.', 'I', 'am', 'not', 'a', 'merry', 'man', '!']
```

```
In [38]: stop_words = set(stopwords.words("english"))
```

```
In [42]: stop_words
```

```
Out[42]: {'a',
          'about',
          'above',
          'after',
          'again',
          'against',
          'ain',
          'all',
          'am',
          'an',
          'and',
          'any',
          'are',
          'aren',
          "aren't",
          'as',
          'at',
          'be',
          'because',
          'been',
          'before',
          'being',
          'below',
          'between',
          'both',
          'but',
          'by',
          'can',
          'couldn',
          "couldn't",
          'd',
          'did',
          'didn',
          "didn't",
          'do',
          'does',
          'doesn',
          "doesn't",
          'doing',
          'don',
          "don't",
          'down',
          'during',
          'each',
          'few',
          'for',
          'from',
          'further',
          'had',
          'hadn',
          "hadn't",
          'has',
          'hasn',
          "hasn't",
          'have',
          'haven',
          "haven't",
          'having',
          'he',
          'her',
          'here',
          'hers',
          'herself',
          'him',
```

'himself',
'his',
'how',
'i',
'if',
'in',
'into',
'is',
'isn',
"isn't",
'it',
"it's",
'its',
'itself',
'just',
'll',
'm',
'ma',
'me',
'mightn',
"mightn't",
'more',
'most',
'mustn',
"mustn't",
'my',
'myself',
'needn',
"needn't",
'no',
'nor',
'not',
'now',
'o',
'of',
'off',
'on',
'once',
'only',
'or',
'other',
'our',
'ours',
'ourselves',
'out',
'over',
'own',
're',
's',
'same',
'shan',
"shan't",
'she',
"she's",
'should',
"should've",
'shouldn',
"shouldn't",
'so',
'some',
'such',
't',
'than',
'that',

```
"that'll",
'the',
'their',
'theirs',
'them',
'themselves',
'then',
'there',
'these',
'they',
'this',
'those',
'through',
'to',
'too',
'under',
'until',
'up',
've',
'very',
'was',
'wasn',
"wasn't",
'we',
'were',
'weren',
"weren't",
'what',
'when',
'where',
'which',
'while',
'who',
'whom',
'why',
'will',
'with',
'won',
"won't",
'wouldn',
"wouldn't",
'y',
'you',
"you'd",
"you'll",
"you're",
"you've",
'your',
'yours',
'yourself',
'yourselves'}
```

```
In [39]: filtered_list = []
```

```
In [40]: for word in words_in_quote:
...     if word.casefold() not in stop_words:
...         filtered_list.append(word)
```

```
In [41]: filtered_list = [word for word in words_in_quote if word.casefold() not in stop_words]
filtered_list
```

```
Out[41]: ['Sir', ',', 'protest', '.', 'merry', 'man', '!']
```

Stemming

```
In [43]: from nltk.stem import PorterStemmer  
from nltk.tokenize import word_tokenize
```

```
In [44]: stemmer = PorterStemmer()  
string_for_stemming = """The crew of the USS Discovery discovered many discoveries.
```

```
In [45]: words = word_tokenize(string_for_stemming)  
words
```

```
Out[45]: ['...',  
'The',  
'crew',  
'of',  
'the',  
'USS',  
'Discovery',  
'discovered',  
'many',  
'discoveries',  
'.',  
'...',  
'Discovering',  
'is',  
'what',  
'explorers',  
'do',  
'.']
```

```
In [46]: stemmed_words = [stemmer.stem(word) for word in words]  
stemmed_words
```

```
Out[46]: ['...',  
'the',  
'crew',  
'of',  
'the',  
'uss',  
'discoveri',  
'discov',  
'mani',  
'discoveri',  
'.',  
'...',  
'discov',  
'is',  
'what',  
'explor',  
'do',  
'.']
```

Tagging Parts of Speech

```
In [60]: from nltk.tokenize import word_tokenize
```

```
In [61]: sagan_quote = """  
... If you wish to make an apple pie from scratch,  
... you must first invent the universe."""
```

```
In [62]: words_in_sagan_quote = word_tokenize(sagan_quote)
```

```
In [26]: nltk.download('averaged_perceptron_tagger')
```

```
[nltk_data] Downloading package averaged_perceptron_tagger to  
[nltk_data] C:\Users\admin\AppData\Roaming\nltk_data...  
[nltk_data] Unzipping taggers\averaged_perceptron_tagger.zip.
```

```
Out[26]: True
```

```
In [63]: nltk.pos_tag(words_in_sagan_quote)
```

```
Out[63]: [('If', 'IN'),  
          ('you', 'PRP'),  
          ('wish', 'VBP'),  
          ('to', 'TO'),  
          ('make', 'VB'),  
          ('an', 'DT'),  
          ('apple', 'NN'),  
          ('pie', 'NN'),  
          ('from', 'IN'),  
          ('scratch', 'NN'),  
          (',', ','),  
          ('you', 'PRP'),  
          ('must', 'MD'),  
          ('first', 'VB'),  
          ('invent', 'VB'),  
          ('the', 'DT'),  
          ('universe', 'NN'),  
          ('.', '.')] 
```

```
In [64]: nltk.download('tagsets')
```

```
[nltk_data] Downloading package tagsets to  
[nltk_data] C:\Users\admin\AppData\Roaming\nltk_data...  
[nltk_data] Unzipping help\tagsets.zip.
```

```
Out[64]: True
```

```
In [65]: nltk.help.upenn_tagset()
```


\$: dollar
 \$ -\$ --\$ A\$ C\$ HK\$ M\$ NZ\$ S\$ U.S.\$ US\$
 '': closing quotation mark
 ' , ' ,
 (: opening parenthesis
 ([{
): closing parenthesis
)] }
 ,: comma
 ,
 -: dash
 --
 .: sentence terminator
 . ! ?
 :: colon or ellipsis
 : ; ...
 CC: conjunction, coordinating
 & 'n and both but either et for less minus neither nor or plus so
 therefore times v. versus vs. whether yet
 CD: numeral, cardinal
 mid-1890 nine-thirty forty-two one-tenth ten million 0.5 one forty-
 seven 1987 twenty '79 zero two 78-degrees eighty-four IX '60s .025
 fifteen 271,124 dozen quintillion DM2,000 ...
 DT: determiner
 all an another any both del each either every half la many much nary
 neither no some such that the them these this those
 EX: existential there
 there
 FW: foreign word
 gemeinschaft hund ich jeux habeas Haementeria Herr K'ang-si vous
 lutihaw alai je jour objets salutaris fille quibusdam pas trop Monte
 terram fiche oui corporis ...
 IN: preposition or conjunction, subordinating
 astride among uppon whether out inside pro despite on by throughout
 below within for towards near behind atop around if like until below
 next into if beside ...
 JJ: adjective or numeral, ordinal
 third ill-mannered pre-war regrettable oiled calamitous first separable
 ectoplasmic battery-powered participatory fourth still-to-be-named
 multilingual multi-disciplinary ...
 JJR: adjective, comparative
 bleaker braver breezier briefer brighter brisker broader bumper busier
 calmer cheaper choosier cleaner clearer closer colder commoner costlier
 cozier creamier crunchier cuter ...
 JJS: adjective, superlative
 calmest cheapest choicest classiest cleanest clearest closest commonest
 corniest costliest crassest creepiest crudest cutest darkest deadliest
 dearest deepest densest dinkiest ...
 LS: list item marker
 A A. B B. C C. D E F First G H I J K One SP-44001 SP-44002 SP-44005
 SP-44007 Second Third Three Two * a b c d first five four one six three
 two
 MD: modal auxiliary
 can cannot could couldn't dare may might must need ought shall should
 shouldn't will would
 NN: noun, common, singular or mass
 common-carrier cabbage knuckle-duster Casino afghan shed thermostat
 investment slide humour falloff slick wind hyena override subhumanity
 machinist ...
 NNP: noun, proper, singular
 Motown Venneboerger Czestochwa Ranzer Conchita Trumplane Christos
 Oceanside Escobar Kreisler Sawyer Cougar Yvette Ervin ODI Darryl CTCA
 Shannon A.K.C. Meltex Liverpool ...
 NNPS: noun, proper, plural

localhost:8888/nbconvert/html/Week11.ipynb?download=false

WDT: WH-determiner
 that what whatever which whichever
 WP: WH-pronoun
 that what whatever whatsoever which who whom whosoever
 WP\$: WH-pronoun, possessive
 whose
 WRB: Wh-adverb
 how however whence whenever where whereby wherever wherein whereof why
 ``: opening quotation mark
 ```

## Lemmatizing

```
In [47]: from nltk.stem import WordNetLemmatizer
```

```
In [50]: nltk.download('wordnet')
```

```
[nltk_data] Downloading package wordnet to

[nltk_data] C:\Users\admin\AppData\Roaming\nltk_data...
```

```
Out[50]: True
```

```
In [52]: nltk.download('omw-1.4')
```

```
[nltk_data] Downloading package omw-1.4 to

[nltk_data] C:\Users\admin\AppData\Roaming\nltk_data...
```

```
Out[52]: True
```

```
In [55]: stemmer.stem('calves')
```

```
Out[55]: 'calv'
```

```
In [56]: lemmatizer = WordNetLemmatizer()

lemmatizer.lemmatize("calves")
```

```
Out[56]: 'calf'
```

```
In []: string_for_lemmatizing = "The friends of DeSoto love scarves."
```

```
In []: words = word_tokenize(string_for_lemmatizing)

words
```

```
In [57]: lemmatized_words = [lemmatizer.lemmatize(word) for word in words]

lemmatized_words
```

```
Out[57]: ['...',
 'The',
 'crew',
 'of',
 'the',
 'USS',
 'Discovery',
 'discovered',
 'many',
 'discovery',
 '.',
 '...',
 'Discovering',
 'is',
 'what',
 'explorer',
 'do',
 '.']
```

```
In [58]: lemmatizer.lemmatize("worst")
```

```
Out[58]: 'worst'
```

```
In [59]: lemmatizer.lemmatize("worst", pos="a")
```

```
Out[59]: 'bad'
```

```
In [67]: lotr_quote = "It's a dangerous business, Frodo, going out your door."
 words_in_lotr_quote = word_tokenize(lotr_quote)
```

```
In [68]: nltk.download("averaged_perceptron_tagger")
 lotr_pos_tags = nltk.pos_tag(words_in_lotr_quote)
 lotr_pos_tags
```

```
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data] C:\Users\admin\AppData\Roaming\nltk_data...
[nltk_data] Package averaged_perceptron_tagger is already up-to-
[nltk_data] date!
```

```
Out[68]: [('It', 'PRP'),
 ("s", 'VBZ'),
 ('a', 'DT'),
 ('dangerous', 'JJ'),
 ('business', 'NN'),
 ('', ''),
 ('Frodo', 'NNP'),
 ('', ''),
 ('going', 'VBG'),
 ('out', 'RP'),
 ('your', 'PRP$'),
 ('door', 'NN'),
 ('.', '.')]

```

## Named Entity Recognition (NER)

```
In [66]: nltk.download("maxent_ne_chunker")
 nltk.download("words")
```

```
[nltk_data] Downloading package maxent_ne_chunker to
[nltk_data] C:\Users\admin\AppData\Roaming\nltk_data...
[nltk_data] Package maxent_ne_chunker is already up-to-date!
[nltk_data] Downloading package words to
[nltk_data] C:\Users\admin\AppData\Roaming\nltk_data...
[nltk_data] Package words is already up-to-date!

NameError Traceback (most recent call last)
Input In [66], in <cell line: 3>()
 1 nltk.download("maxent_ne_chunker")
 2 nltk.download("words")
----> 3 tree = nltk.ne_chunk(lotr_pos_tags)

NameError: name 'lotr_pos_tags' is not defined
```

```
In [69]: tree = nltk.ne_chunk(lotr_pos_tags)
```

```
In [70]: tree.draw()
```

```
In [71]: tree = nltk.ne_chunk(lotr_pos_tags, binary=True)
tree.draw()
```

```
In [72]: quote = """
... Men like Schiaparelli watched the red planet—it is odd, by-the-by, that
... for countless centuries Mars has been the star of war—but failed to
... interpret the fluctuating appearances of the markings they mapped so well.
... All that time the Martians must have been getting ready.
...
... During the opposition of 1894 a great light was seen on the illuminated
... part of the disk, first at the Lick Observatory, then by Perrotin of Nice,
... and then by other observers. English readers heard of it first in the
... issue of Nature dated August 2."""
```

```
In [75]: def extract_ne(quote):
... words = word_tokenize(quote)#, Language=Language)
... tags = nltk.pos_tag(words)
... tree = nltk.ne_chunk(tags, binary=True)
... return set(
... " ".join(i[0] for i in t)
... for t in tree
... if hasattr(t, "label") and t.label() == "NE"
...)
```

```
In [76]: extract_ne(quote)
```

```
Out[76]: {'Lick Observatory', 'Mars', 'Nature', 'Perrotin', 'Schiaparelli'}
```

```
In []:
```